# QS101: Introduction to Quantitative Methods in Social Science

## Week 15: Measures of Association: Correlation

### Dr. Florian Reiche

Teaching Fellow in Quantitative Methods
Course Director BA Politics and Sociology
Deputy Director of Student Experience and Progression, PAIS

February 5, 2015

Recap

Correlation

Dr. Florian Reiche

QS101: Introduction to Quantitative Methods in Social Science

Recap

# Queries

- What is $\chi^2$?

# Queries

- What is $\chi^2$?
- What are degrees of freedom?

## Queries

- What is $\chi^2$?
- What are degrees of freedom?
- What is the p-value?

# Numbers and the Media

https://www.youtube.com/watch?v=oDPCmmZifE8

Correlation

## Definition

- Correlation is a statistical tool that determines the degree of relationship between two different variables
- If correlation is strong, a person's score on one variable helps us predict the person's score on another variable
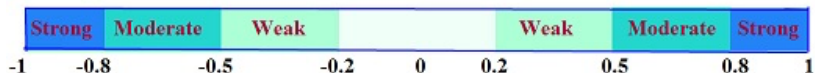- It is limited with the range of $-1$ and $+1$

# Classification of the Relationship

**Pearson Correlation Coefficient**

| Negative Correlation | No Correlation | Positive Correlation |
|---|---|---|

| Strong | Moderate | Weak | | Weak | Moderate | Strong |
|---|---|---|---|---|---|---|

-1      -0.8           -0.5              -0.2      0      0.2           0.5              0.8      1

Dr. Florian Reiche

QS101: Introduction to Quantitative Methods in Social Science

## Strong Positive Relationship

- ▶ The higher the score on one variable, the higher the score on the other variable
- ▶ The lower the score on one variable, the lower the score on the other variable.

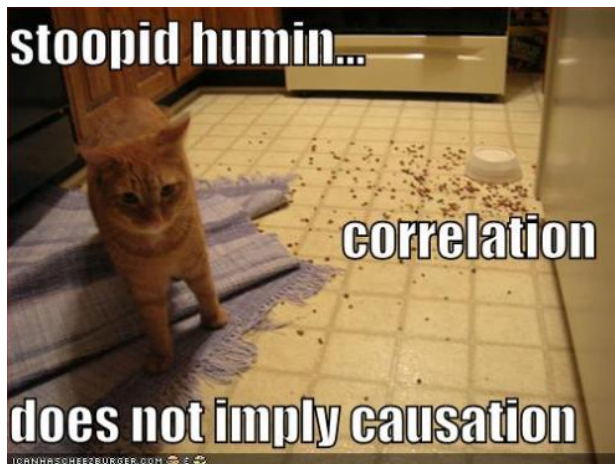Example: Time spent on revision, and exam mark.

## Strong Negative Relationship

- The higher the score on one variable, the lower the score on the other variable
- The lower the score on one variable, the higher the score on the other variable.

Example: Time spent in the Dirty Duck, and module mark.

# No Relationship

- Correlation coefficient (r) = 0
- Here, the score on one variable tells you nothing about the score on the other variable

## Correlation and Causation

# Example 1

# Example 2

# Example 3

## Why bother then?

- You can never show a cause-effect relationship with correlation
- Yet, you can get an idea about the data, and patterns within in
- This can help you develop ideas about cause-effect relationships

## 4 Types of correlation

- ▶ Pearson's $r$: A measure of the strength of a relationship between two continuous variables
- ▶ Spearman's $r$: A measure of the similarity between two ordinal rankings of a single set of data
- ▶ Point-biseral $r$: A measure of the strength of the relationship between one continuous variable and one dichotomous variable (e.g. gender, democracy, etc.)
- ▶ Phi ($\phi$) correlation: A measure of the strength of the relationship between two dichotomous variables

## The Pearson Product-Moment Correlation Coefficient

Is the most commonly employed measure:

$$r = \frac{N\Sigma xv - (\Sigma x)(\Sigma y)}{\sqrt{(N\Sigma x^2 - (\Sigma x)^2)(N\Sigma y^2 - (\Sigma y)^2)}} \tag{1}$$

N: Number of pairs of scores

## Example

| Subject | Cigarettes | Years Lived |
|:-------:|:----------:|:-----------:|
| 1 | 25 | 63 |
| 2 | 35 | 68 |
| 3 | 10 | 72 |
| 4 | 40 | 62 |
| 5 | 85 | 65 |
| 6 | 75 | 46 |
| 7 | 60 | 51 |
| 8 | 45 | 60 |
| 9 | 50 | 55 |

# Example contd.

- For example:
  $\Sigma x = 25 + 35 + 10 + 40 + 85 + 75 + 60 + 45 + 50 = 425$

- etc.

$$r = \frac{(9)(24,640) - (425)(542)}{\sqrt{((9)(24,525) - (425)^2)((9)(33,188) - (542)^2)}}$$

$$r = -0.6111$$

$$r = -0.61$$

# The Pearson Product-Moment Correlation Coefficient

- ▶ Obtained for a sample drawn from the population, denoted $r$
- ▶ The population value is called rho ($\rho$)
- ▶ We are therefore interested in:
    - ▶ $H_0 : \rho = 0$
    - ▶ $H_a : \rho \neq 0$

## Significance Test

- ▶ Uses the t-distribution (if you do not know what this is, read up on it, now!)
- ▶ The t-test formula for a correlation coefficient is as follows:

$$t = \frac{r}{\sqrt{\frac{1-r^2}{N-2}}} \qquad (2)$$

- ▶ This would give you the critical value (just like with $\chi^2$ last week)
- ▶ Calculate the degrees of freedom (here: $df = N - 2$)
- ▶ You choose a level of significance, find the critical value, compare the values and decide

## Example contd.

$$t = \frac{r}{\sqrt{\frac{1-r^2}{N-2}}}$$

$$t = \frac{-0.6111}{\sqrt{\frac{1-(-0.6111)^2}{9-2}}}$$

$$t = -2.042$$

- We want 95% significance level
- Critical value for $df = 7$: $t = +2.365$ and $-2.365$
- Is this significant?

## Example contd.

$$t = \frac{r}{\sqrt{\frac{1-r^2}{N-2}}}$$

$$t = \frac{-0.6111}{\sqrt{\frac{1-(-0.6111)^2}{9-2}}}$$

$$t = -2.042$$

- ▶ We want 95% significance level
- ▶ Critical value for $df = 7$: $t = +2.365$ and $t = -2.365$
- ▶ Is this significant? No!

## The other 3 Measures

- Spearman: Coolidge, pp. 204-206
- Point Biseral: Coolidge, pp. 208-212
- $\phi$ correlation: Coolidge, pp. 212-213